

01TXFSM - Machine Learning and Deep Learning

Final Project First Person Action Recognition

Eros Fani - s269781
Politecnico di Torino

eros.fani@studenti.polito.it

Gabriele Trivigno - s276807
Politecnico di Torino

gabriele.trivigno@studenti.polito.it

Cristiano Gerbino - s277058
Politecnico di Torino

s277058@studenti.polito.it

Abstract

1. Introduction

The main scope of this project was the field of First Person Action Recognition, which is one of the up-and-coming domains in modern Computer Vision, due to the recent spread of wearable devices, often camera-equipped. Even though action recognition and video analysis in general have been topics of research for quite some time now, this particular task introduces some newly faced challenges such as the strong egomotion component, which is inevitable as cameras are placed directly on the actor. Moreover the First Person point of view lacks some important information about the actor and its pose, showing only a part of the scenario, making inference harder. The most frequently adopted techniques to tackle this issues combine spatial information, taken from RGB frames, and motion clues, extracted in various ways, such as temporal stream networks based on optical flow, attention modules or 3D CNNs. We will discuss in depth some of this methodologies, highlighting in our opinion what are the weaknesses of such approaches, and proposing possible ways to circumvent them.

1.1. Goals

The first goal of the project was to replicate the main results of [?]. Afterwards, we integrated the cited architecture with the self-supervised block proposed in [?], implementing it in 2 different ways.

After trying to replicate these experiments with the hyperparameters setting reported in [?], with the experience acquired we chose other sets of hyperparameters that could potentially fit better our dataset of interest, and so performed a grid search upon them.

At last we have tried to improve the performances of the results of [?] and [?] with some innovative ideas of our own.

1.2. Our contribution

2. Related works

Existing literature has underlined how the most crucial information to encode into features in order to obtain good performances on the task of First Person Action Recognition are the hands motion of the actor, and the appearance of objects manipulated, as well as their interaction. Putting together this

2 fundamental pieces of information is perhaps the most challenging task, and more importantly doing so without complicating too much the architecture, which could lead to high computational costs, cutting of potential on-line applications. [?] proposes a spatial attention mechanism to focus on the region consisting of the object, arguing that objects handled are strongly representative of egocentric activities. To keep track of said objects throughout frames, this spatial attention is paired with a convLSTM module, whose purpose is to temporal encode the frame-level features formerly extracted. In addition to this, they also use a temporal network fed with optical flow images, following the purpose of better encoding motion changes in the considered videos. The task is then addressed with a late fusion of the 2 outputs with a dense layer to finally perform the classification step. Whilst reasonably successful, this approach consists of a two stream network, which requires several training stages resulting in a massive number of parameters to train. [?] argues that adding a self-supervised block which leverages motion maps at the output of the CNN backbone can provide the convLSTM module with a richer input that already encloses motion information. They show that in this way the network is able to obtain a meaningful enough representation of both appearance and movement to achieve state-of-the-art results, so that the temporal network is no longer needed. Many other works recognize the importance of slimming the architecture to avoid having to train a 2 stream network. [?] introduces a unified model that aims at representing spatio-temporal relationship using only RGB frames, eliminating the heavy computation time required to extract optical flow, which makes the most difference at test time, once the model has been deployed. In their MFNet they insert motion filter blocks in the middle of CNN layers, that act upon feature maps extracted from shared-networks feed-forwarded by two consecutive input frames. The flow is estimated moving each channel of these feature maps in a different spatial direction $\delta := (\Delta x, \Delta y)$. [?] argues that the procedure described, though advantageous parameter-wise and computation-wise, suffer from inferior performance compared with more powerful two-stream networks. They maintain such performance gap is due to the iterative optimization that standard optical flow methods implement and that [?] fails to reproduce. To overcome this issue they propose a new, fully-differentiable, layer, built to extract flow from any CNN feature map. By learning the flow parameters of a smaller resolution CNN tensor, in an end-to-end fashion together with the CNN training, they claim to achieve

the same representational power as traditional flow methods, and even better since flows optimized for activity recognition is different from true optical flow. This last two works introduce a great way to encode motion clues in the features extracted, but none of them combine this information with a temporal representation of those features. Whereas [?] fails to exploit RGB frames to get motion informations beyond appearances. Even though [?] tries to overcome this with the motion segmentation task, we found that their approach can sometimes lacks of precision in identifying the region of interest, since it acts on very high-level features. In this project we propose three new different approaches. First of all, we use an improved version of the self-supervised block adopted in [?] to integrate in the work of [?], eliminating the need of a second network and improving the granularity of the motion encoding using higher resolution maps matched with lower level features. Afterwards we adopt the representation-flow layer proposed by [?] together with a temporal encoding module to better exploit temporal relationships in the features extracted. Finally we experimented with a different approach: encoding the concept of dynamics directly into the convLSTM module, adding a discriminator classifier that has been trained to discern whether it has been fed with static frames, or actual video frames.

2.1. Dataset description

The dataset under analysis is a modified version of GTEA61¹. The dataset contains the videos in form of frames, and also two kind of preprocessed images: *motion maps* and *optical flows*. The folder schema of the dataset is shown in Figure ???. Videos represent 61 class actions performed by 4 different users (*S1, S2, S3, S4*). Sometimes for some actions more than one video is available. The total number of videos in the dataset is, however, 457, which actually means that it is a quite small dataset.

The optical flow methods try to calculate the motion between two image frames which are taken at times t and $t + \Delta t$ at every voxel position. The warp flow methods try also to remove the motion of the wearable camera. We have two kind of these last representations in our dataset: one computed in the horizontal axis (folder *flow.x.processed*) and one other computed in the vertical axis (folder *flow.y.processed*).

¹Georgia Tech Egocentric Activity Datasets: <http://cbs.ic.gatech.edu/fpv/>

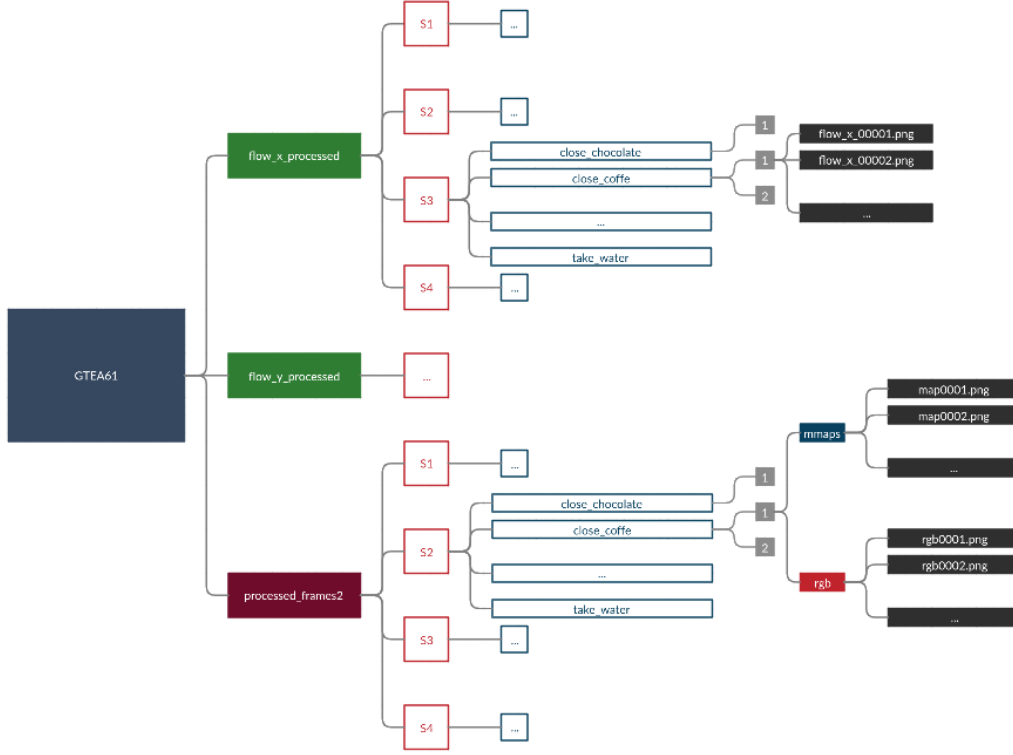


Figure 1: Folder schema of our GTEA61

The motion maps are special black-and-white images which represent the spatial location in which the Motion Segmentation task of [?] focuses its attention per each frame. The mmaps present large similarities with the warp flows.

The differences between the kind of available images in our dataset are shown in Figure ??.

2.2. Data cleaning

The dataset was almost clean already from the beginning, but we encountered two problems within it:

- there were hidden useless folders *.DSstore* inside each one of the user folders. These have been removed
- some of the first mmaps of some videos were missing. In these cases we have simply duplicated the second mmap

3. Descriptions of the models

Here we describe the models that we have used to perform our experiments.

3.1. Egornn

Egornn is a Recurrent Neural Network. The overall architecture of *Egornn* is shown in Figure ??. This net is based on *resnet34* [?], which constitutes the main block. *resnet34* has five convolutional layers inside itself: with respect to Figure ?? they are: *Conv*, *Layer1*, *Layer2*, *Layer3* and , *Layer4*. From now on we will refer to these blocks respectively *conv1*, *conv2*, *conv3*, *conv4* and *conv5*.

At the termination of the *resnet34* is placed a *Spatial Attention Layer*. It calculates a *Class Activation Map* (CAM) that is capable to identify the image regions that have been used by the CNN to identify

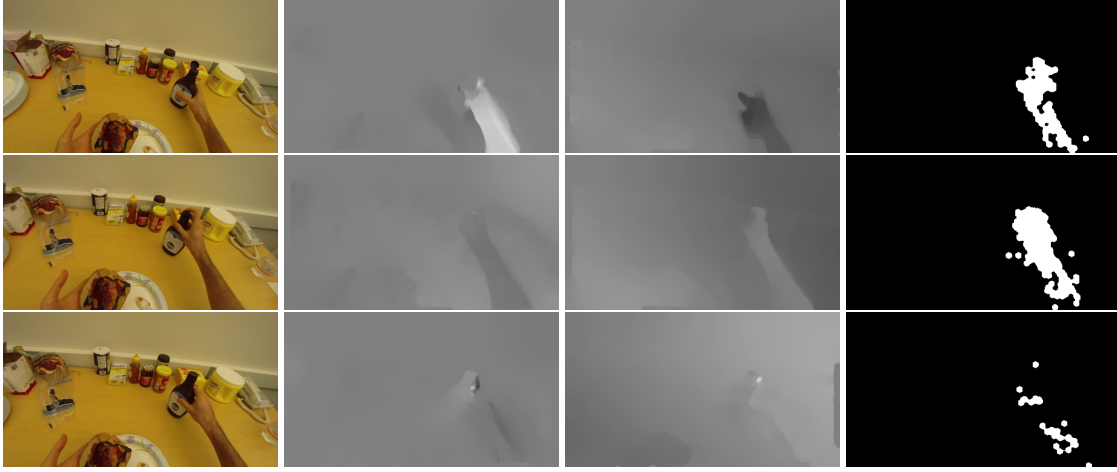


Figure 2: Types of images in our dataset. In this example is shown a sample of images from the *close_chocolate* action. From the left column to the right column: rgbs, warp flows x, warp flows y, motion maps

the class under analysis. It is computed by taking the output of the *softmax* layer and the output of *conv5* and taking the linear combination of all the weights of *conv5* and the weights of the softmax.

The output of the CAM is then sent to a *softmax* layer to obtain a probability map, which is called *Spatial Attention Layer* (SAM). The output of the SAM is finally multiplied, cell by cell (Hadamard product), with the output of *conv5*, obtaining another tensor of weights which is sent to a *Convolutional Long Term Support Memory* block (ConvLSTM).

The reason for the usage of the ConvLSTM block is that, up to now, what the net does is to take each frame and to try to make predictions based only on the features that the net can extract from those frames, without taking into consideration the temporal encoding of frame level features. The convLSTM block take into consideration, for each frame i , both the output of the SAM for the layer i and the output of the ConvLSTM for the layer $i - 1$, constituting a recursive structure.

The last output of the ConvLSTM (the output obtained from the last frame of a particular video) is average pooled and reshaped to obtain a final classification layer with 61 neurons (i.e. the number of classes of our dataset).

3.2. Flow_resnet34

Flow_resnet34 is just a *resnet34* edited to work with the warp flows. It gets five warp flows from *processed_frames_x* and five from *processed_frames_y* in form of a tensor of ten channels and tries to make predictions on the 61 classes.

3.3. Two stream model

Egornn learns appearance features, while *flow_resnet34* learns motion features. The way to join the two nets is to concatenate the two output layers and to add at the end a fully connected layer to get the class category scores.

3.4. Motion Segmentation branch applied to egornn

The problem which [?] tries to overcome is that in the two stream model motion and appearance are actually separately learned, without taking into account the spatial-temporal relationships.

We have built an architecture similar to *sparnet*, where the *motion segmentation block* is the same but the *action recognition block* has been substituted by *egornn* (like in one of the attempts in [?]). The architecture is shown in Figure ???. We have

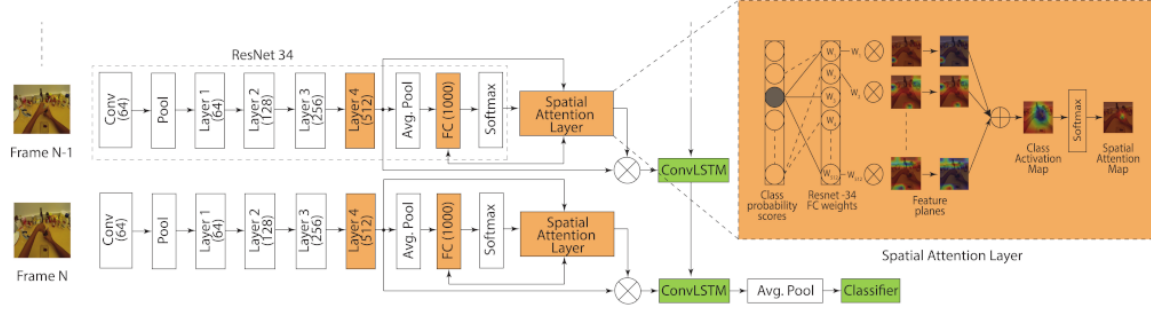


Figure 3: Architecture of *egornn*

used this architecture with some granular variations during our experiments, but the main blocks are always as shown in Figure ?? . The input of the convolutional layer of MS Block is taken from one of the convolutional layers of *resnet34* of *egornn* (the actual layer varies with our experiments). Then, after the convolutional layer, there is a fully connected layer followed by a softmax which normalizes the weights between 0 and 1. *mmaps ground truth* and *rgb ground truth* represent the mmaps and the rgb after the transformations. The transformations applied to the mmaps are the same applied to the rgbs, plus a small amount of proper mmaps transformations which always ends with a transformation which linearizes the pixels (from a 2 dimensional tensor per mmap we get a 1 dimensional tensor per mmap). For the msblock these linearized pixels represent a real ground truth, because each of the output neurons of the MS Block is used to predict the values of the mmaps ground truth. The pixel losses are summed together (obtaining as result L_{ms}) and then are summed again with the *egornn* loss (L_c). The final loss is used to compute the gradients to update the weights.

3.5. Static-dynamic discriminator

Starting from the model described above, we have added a final binary classifier to *egornn* after the convLSTM, parallel to the other classifier already present which still tries to predict the actual class of the video. The idea was to force the net to learn the motion features from the rgbs. During the training phase this classifier gets 2 kinds of sequences of frames: one is the same of the original classifier, while the other gets a sequence of identical frames. This discriminator should be able to recognize the

actual videos from the static frames. In this way the gradients should adapt to focus the attention on the motion.

4. Experiments

Our nets are always trained on a predefined train set, which includes all and only the videos of the users *S1*, *S3* and *S4*, while validation and test sets coincide and is constituted by all and only the videos of a single user, *S2*. In addition, the weights of the *resnet34* are pretrained on ImageNet. Each model is always validated while it is trained, so for each training phase we selected the weights with the highest accuracy at a particular epoch as the best ones.

Due to Colab limitations of GPU memory, we have only been able to perform experiments on a limited amount of frames (7 or, in less cases, 16). Due to this problem our results should be interpreted not as absolute value of the accuracy, but as a sort of relative value with respect to the number of frames for each video in our batches.

The size of our batches has always been left to 32, as well as the number of hidden units of the convLSTM module, fixed at 512. Our optimization algorithm is always Adaptive moment estimation (ADAM) with the only exception of *flow_resnet34*, for which it is Stochastic Gradient Descent (SGD). When using this last optimizer, the momentum has always been left to 0.9. The scheduler is a MultiStepLR scheduler, which decreases the original learning rate LR by a factor GAMMA at each value of STEP_SIZE.

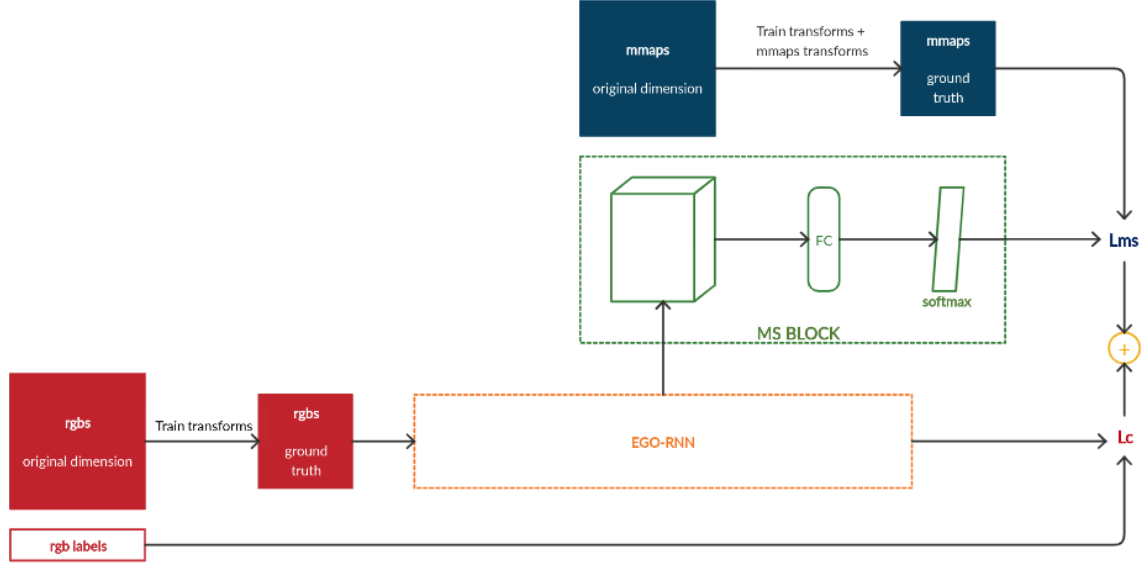


Figure 4: Generic architecture of motion segmentation branch applied to *egornn*

4.1. Egornn

We have replied some of the same experiments of [?] on the original *egornn*. We have run each of these experiments three times and then we have averaged the results.

First, we have performed the classification by using the *egornn* without and with the CAM. The training phase has been divided in two parts, as in the original paper:

1. train of ConvLSTM and Classifier (green blocks in Figure ??)
2. train of conv5 (layer4 of *resnet34*), FC(1000), Spatial Attention Layer (orange blocks in Figure ??) in addition to the previously listed blocks

The values of the hyperparameters for the first stage are:

LR	1e-3
WEIGHT_DECAY	4e-5
NUM_EPOCHS	200
STEP_SIZE	[25, 75, 150]
GAMMA	0.1

While, for the second stage, they are:

LR	1e-4
WEIGHT_DECAY	4e-5
NUM_EPOCHS	150
STEP_SIZE	[25, 75]
GAMMA	0.1

Then, we have also trained *flow_resnet34* alone. In this case we used only 5 frames per each flow (x and y) due to the fact that for some videos no more than 5 frames were provided.

The values of the hyperparameters in this case are:

LR	1e-2
WEIGHT_DECAY	5e-4
NUM_EPOCHS	750
STEP_SIZE	[150, 300, 500]
GAMMA	0.5

At last we performed the two stream training with the following values for the hyperparameters:

Configurations	Frames	Mean accuracy
EGO-RNN without CAM - stage 1	7	29.89
EGO-RNN without CAM - stage 1	16	27.87
EGO-RNN without CAM - stage 2	7	50.00
EGO-RNN without CAM - stage 2	16	50.57
EGO-RNN - stage 1	7	41.38
EGO-RNN - stage 1	16	46.84
EGO-RNN - stage 2	7	58.91
EGO-RNN - stage 2	16	65.52
flow_resnet34	5	46.26
two-stream (joint train)	7*	57.76
two-stream (joint train)	16*	66.38

Figure 5: Summary of the results over different configurations. Each value of the mean accuracy is the mean of the accuracies over three identical experiments. *the number of frames refers to the *egornn* branch (for the *flow_resnet34* branch the number of frames is always 5)

LR	1e-2
LR_FLOW	1e-4
WEIGHT_DECAY	5e-4
NUM_EPOCHS	250
STEP_SIZE	[1]
GAMMA	0.99

Where LR is the learning rate of *egornn* and LR_FLOW is the learning rate of *flow_resnet34*.

The summary of our results is shown in Figure ?? . From here it raises that the best model is the two-stream (joint train) with 16 frames, followed by EGO-RNN - stage 2 with 16 frames, which is behind the two-stream model by less than 1 point of mean accuracy. Due to the averaging between three identical runs we can rely on this result and assert that the contribution of *flow_resnet34* slightly increases the performances, but also that the most of the contribution is given by *egornn*.

In Figure ?? and Figure ?? are shown respectively the validation accuracy and the validation loss by epoch of one random extracted run per each one of the attempts with 16 frames (5 in case of *flow_resnet34*) and only for the stage 2 when a two stage training is required.

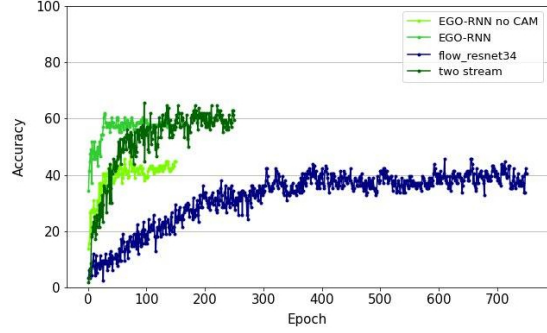


Figure 6: Validation accuracy by epoch of one random extracted run for the four most interesting training configurations

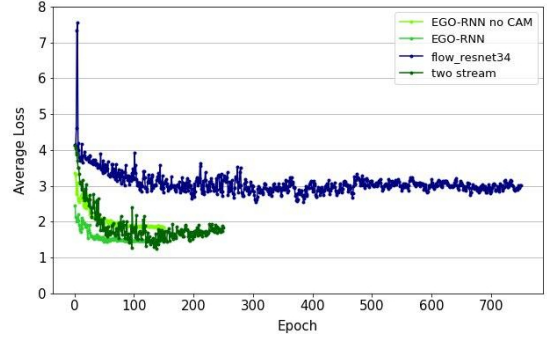


Figure 7: Validation loss by epoch of one random extracted run for the four most interesting training configurations. The losses are the average of each batch loss within a single epoch

From Figure ?? and Figure ?? is even more evident that *flow_resnet34* is highly inefficient alone, and that the results with the CAM are heavily better than the results without the CAM (higher accuracy and lower loss at every epoch). The two-stream model requires more time to get high accuracies, and overall it seems to have the same behaviour of *egornn* when at full capacity, but it is noisier and so it is easier that for some epoch the accuracy is higher.

Stage	Frames	Mean accuracy
1	7	44.83
1	16	50.86
2	7	60.63
2	16	62.07

Figure 8: Mean accuracies over three identical experiments with the same fixed set of values for the hyperparameters, by varying only the stage of training and the number of frames - classification experiment

4.2. Motion Segmentation branch applied to egornn

First of all we have replied the same experiment of [?] (the one in which *egornn* is the Action Recognition Block). The hyperparameters used are:

LR	1e-3
WEIGHT_DECAY	4e-5
NUM_EPOCHS	150
STEP_SIZE	[50, 100]
GAMMA	0.1

for the stage1, while they are the following:

LR	1e-4
WEIGHT_DECAY	4e-5
NUM_EPOCHS	150
STEP_SIZE	[25, 75]
GAMMA	0.1

for the stage2. We have decreased the number of epochs for the stage1 because we have observed that going too far with the epochs, with the loss decreased after the steps of the optimizer, would be meaningless and does not give any significant result.

The results with the same values for the hyperparameters are shown in Figure ??.

As expected, the best performances are achieved after the second stage of training and with 16 frames.

We have also replied this experiment as a regression problem. [REGRESSION IMPLEMENTATION]

Stage	Frames	Mean accuracy
1	7	44.83
1	16	50.86
2	7	59.77
2	16	66.38

Figure 9: Mean accuracies over three identical experiments with the same fixed set of values for the hyperparameters, by varying only the stage of training and the number of frames - regression experiment

In Figure ?? are shown the results.

Also in this case the highest value of the accuracy is obtained for 16 frames. The first stage is exactly the same so we have not replied it three times more. For the second stage we observe that with this combination of values for the hyperparameters the regression performs more than 4 points better than classification with 16 frames, while the accuracies are roughly the same when the number of frames is 7.

We have performed also a complete grid search to improve the performances of the net.

Before starting, it has been executed a brief tuning step for the stage 1 of training, which is the same for both the classification and the regression methodologies. The values of interest have been STEP_SIZE and GAMMA only. From this we obtained that the optimal values for the stage1 are exactly the same used in the previous experiment.

After this preliminary step, the complete grid search has been executed only on the stage 2, which is the most important one. The number of selected frames is 16, for which we have seen that the performances are better than for 7. We have reduced NUM_EPOCHS to 100 to save computational time (in addition the highest accuracy is rarely obtained after epoch 100). For the classification strategy the results are shown in Figure ??, while for the regression they are shown in Figure ??.

As we can observe the best performances are achieved with:

- $LR = 10^{-4}$, $WEIGHT_DECAY = 4 \cdot$

LR	WEIGHT_DECAY	GAMMA	STEP_SIZE	0.1	0.2	0.5
1e-04	4e-03		[30, 70]	73.28	64.66	65.52
1e-04	4e-03		[40, 80]	67.24	75.00	67.24
1e-04	4e-03		[50, 100]	67.24	67.24	68.10
1e-04	4e-05		[30, 70]	62.07	68.97	65.52
1e-04	4e-05		[40, 80]	67.24	68.97	66.38
1e-04	4e-05		[50, 100]	65.52	67.24	67.24
5e-04	4e-03		[30, 70]	68.10	64.66	70.69
5e-04	4e-03		[40, 80]	65.52	62.93	63.79
5e-04	4e-03		[50, 100]	62.45	65.52	66.38
5e-04	4e-05		[30, 70]	64.66	65.52	61.23
5e-04	4e-05		[40, 80]	65.52	66.38	62.07
5e-04	4e-05		[50, 100]	68.97	62.45	67.24

Figure 10: Accuracies at various combinations of hyperparameters for the classification strategy

LR	WEIGHT_DECAY	GAMMA	STEP_SIZE	0.1	0.2	0.5
1e-04	4e-03		[30, 80]	62.07	62.93	62.07
1e-04	4e-03		[40, 90]	66.38	68.97	64.66
1e-04	4e-05		[30, 80]	60.34	62.07	67.24
1e-04	4e-05		[40, 90]	62.93	66.38	64.66
5e-04	4e-03		[30, 80]	64.66	61.23	67.24
5e-04	4e-03		[40, 90]	65.52	66.38	71.55
5e-04	4e-05		[30, 80]	68.10	56.03	72.41
5e-04	4e-05		[40, 90]	66.38	65.52	68.10
5e-05	4e-03		[30, 80]	61.21	63.79	62.07
5e-05	4e-05		[30, 80]	61.21	62.07	66.38

Figure 11: Accuracies at various combinations of hyperparameters for the regression strategy

10^{-3} , $STEP_SIZE = [40, 80]$ for the classification method, for an accuracy of 75.00

- $LR = 5 \cdot 10^{-4}$, $WEIGHT_DECAY = 4 \cdot 10^{-5}$, $STEP_SIZE = [30, 80]$ for the regression method, for an accuracy of 72.41

In general we can say that the performances are better with the classification method, and this is a surprisingly result based on the observations with-

out the grid. The reason for the previous result was that, as emerges from Figure ?? and Figure ??, the optimal values for the hyperparameters falls in a total different region with respect to the two different methodologies.

In Figure ?? we can observe the effects of the values of the hyperparameters on the accuracy and on the loss of both train and validation and over three different experiments present in the grid search with the classification methodology. The colors of the graphs underline the different learning rate before and after each step of the optimizer. The steps are shown on the x axis and marked with a vertical gray line. In (a) is shown the best result obtained, while (b) and (c) are representative examples with different tuned values for the hyperparameters. In (a) the loss already starts from a low value and the accuracy already starts from a high value. LR highly influences the starting point of the accuracy and the loss in epoch 1 (higher LR \rightarrow lower initial accuracy and higher initial loss). To show this behaviour it has been picked the example in (b), which has the same STEP_SIZE but different values for the other hyperparameters. We can observe that, from Figure ??, there is not a tangible difference between the column with GAMMA = 0.1 and the column with GAMMA = 0.2, while it is more evident that with GAMMA = 0.5 the accuracy worsen a bit. So, it has been selected another example, (c), with GAMMA = 0.1 (which is a different gamma with respect to (a)) and the same values of LR and WEIGHT_DECAY (which differs from (a)), and with another STEP_SIZE. The effect of this slight variation on the STEP_SIZE seems to be not so meaningful, while the actual main characters of the different behaviour of the net are LR and WEIGHT_DECAY, together with GAMMA which is good when it is 0.1 as well as when it is 0.2. As last observations:

- is evident from all the three experiments that the train accuracy continue growing also when the validation accuracy starts flattening, and the behaviour is the relative one for the loss
- both the losses and the accuracies are “noisier” when the learning rate is higher
- the global behaviour is strongly influenced by the random batches which feed the net, and so the results should be intended within an error margin

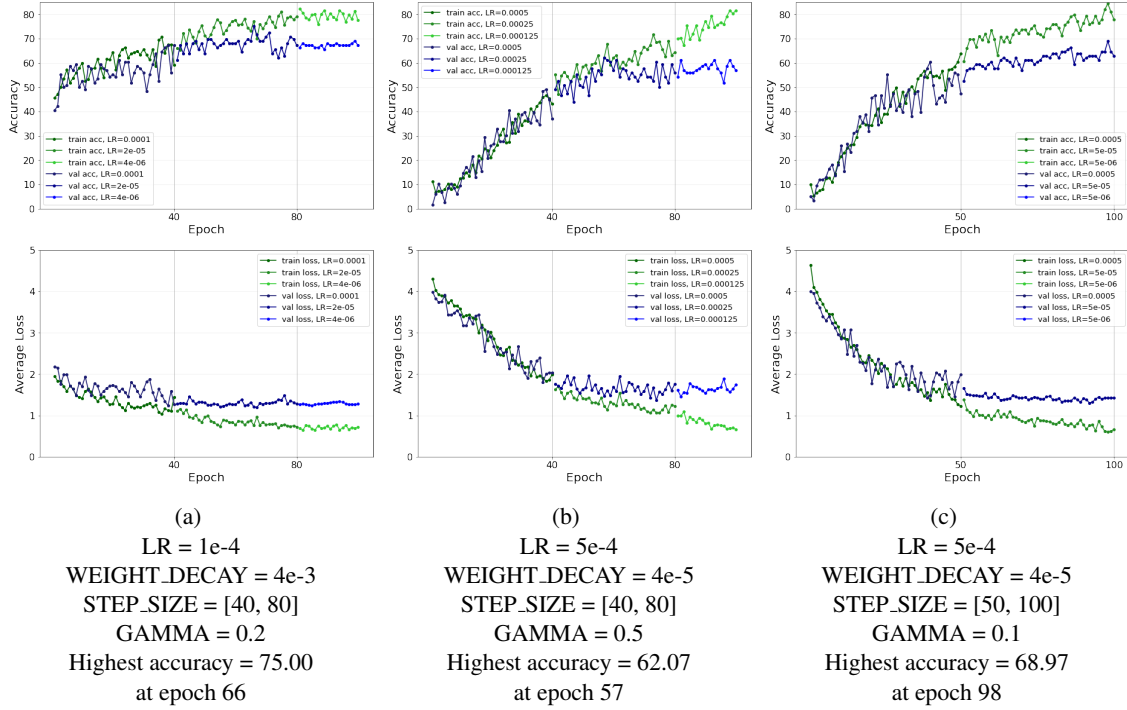


Figure 12: Comparison between three representative experiments in the grid search with the classification method

4.3. New experiments

4.3.1 The problem of the Motion Maps

To assert if there is an actual contribution of the Motion Segmentation branch and how much it is significant for the final accuracy we analyzed the loss on the mmaps with a special attention. An example of this loss is shown in Figure ?? . As we can see, the loss on the train set stops to slightly drop after a bunch of epochs.

We better studied this behaviour by analyzing the CAMs.

To overcome these problems we have tried to increase the resolution of the downsampled mmaps and to treat the problem as.

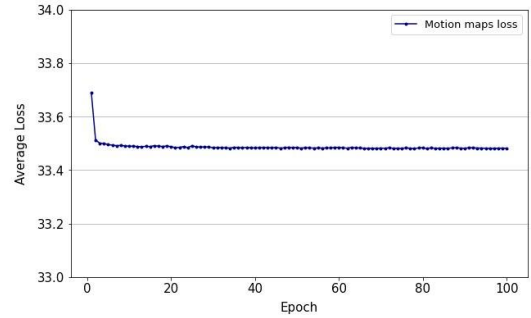


Figure 13: An example of the loss on the classification of the pixels of the downsampled mmaps performed by the Motion Segmentation Branch

4.3.2 Static-dynamic discriminator

The reasons for this approach are explained in Section ?? . Also here, due to *Colab out of memory*

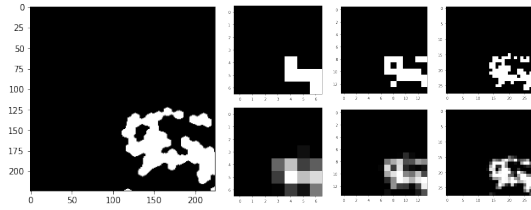


Figure 14

- [5] S. Sudhakaran and O. Lanz. Attention is all we need: Nailing down object-centric attention for egocentric activity recognition, 2018.

errors, we had to perform the experiments with 7 frames per video instead than with 16. We have tried four different configurations for the training pipeline, each one tried three times to give consistency to the results:

1. training of *conv5*, *convLSTM* and the final static-dynamic classifier when forwarding the static frames
2. training of *convLSTM* and the final static-dynamic classifier when forwarding the static frames
3. training *convLSTM* and the final static-dynamic classifier when forwarding the static frames. The static-dynamic loss is multiplied by a constant $\text{ALPHA} = 0.5$
4. training *convLSTM* and the final static-dynamic classifier when forwarding the static frames. The static-dynamic loss is multiplied by a constant $\text{ALPHA} = 0.1$

4.3.3 Flow di Gabriele

5. Problems occurred and possible ideas

6. Conclusion

References

- [1] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2015.
- [2] M. Lee, S. Lee, S. Son, G. Park, and N. Kwak. Motion feature network: Fixed motion filter for action recognition, 2018.
- [3] A. Piergiovanni and M. S. Ryoo. Representation flow for action recognition, 2019.
- [4] M. Planamente, A. Bottino, and B. Caputo. Joint encoding of appearance and motion features with self-supervision for first person action recognition, 2020.